# Run II Computing Review
June 4-6 2002

## *Executive Summary*

The review of Fermilab Run II Computing was held at Fermilab on June 4-6, 2002. The members of the review committee were:

- Ian Bird, Jefferson Lab (chair)
- Elizabeth Buckley-Geer, FNAL
- John Hobbs, Stony-Brook
- Richard Mount, SLAC
- Pekka Sinervo, University of Toronto

The charge to the committee was to review the plans of CDF and D0 for providing the capabilities to collect, store, and analyse the Run II data, and to assess whether the tentative budgets for equipment and tapes are adequate. The experiments had been given guidelines of $2M for equipment and $500K for tapes per year per experiment.

The review team heard presentations from members of both CDF and D0 as well as representatives from the Computing Division.

Both experiments have their computing systems in production operation and are taking and analyzing data. The problems with the CDF data handling system have been resolved by moving to Enstore – this was possible only with the support of the CD/ISD group and good collaboration, and the transition to a stable working system was achieved in only 6 months.

Rather than write separate reports on each experiment's plans, we chose to look at several broad areas which are reflected in the structure of the remainder of this document: Data Handling, Farms and Production facilities, Central Analysis Facilities, Remote Analysis plans including Grid, Infrastructure, and Management (Budgets, Staffing, Schedules, Planning).

## Summary of Recommendations

- Proceed for now on the basis of the plans presented here, but maintain sufficient flexibility over the next year as needs and assumptions become better understood so that an better optimized system can be developed over the period of Run II.
- CDF continue the deployment of commodity based fileservers, and if proven to work, it is recommended that DZero consider this as a replacement of the SMP system.
- The experiments and the Computing Division consider all costs, including maintenance and personnel, when comparing various system options.
- Both collaborations develop more detailed plans for the coordinated use of remote computing facilities, and identify how best they can be used to meet the overall Run II data processing and physics analysis needs. We recommend that the collaborations make this contribution more formal, perhaps in the form of MOU's.
- The Laboratory work with ESNET to develop a high-speed connection to STARLIGHT on the timescale of the next 6 months in support of remote analysis capabilities.

- CDF should be encouraged to contribute dedicated effort to the SAM project, since now both experiments are making use of it.
- The commodity-based analysis cluster for D0 should be a centrally managed system.

## *Data Handling*

Fermilab now has sound foundations for Run II data handling. The Enstore mass storage software can be confidently expected to meet the Run II needs and both experiments have a viable baseline model for tape storage and retrieval based on the use of the ongoing series of high-density tape drives from STK. Options to use the cheaper LTO drives will be considered. At the file-handling level, D0 is firmly committed to the SAM system developed jointly with CD, and CDF is evaluating SAM. Although not yet tested at a stress level that will be typical of Run-2 production and analysis, SAM appears to be on track to achieve the required performance and robustness. We believe that use of SAM by both experiments is clearly the right choice. The acronym SAM means "Sequential Access via Metadata". However, SAM itself handles only files and does not constrain the way in which these files will be accessed once they have been made available. It was notable that, in spite of this, sequential access within the CDF PAD and D0 DST datasets is still the baseline approach of both experiments.

At levels above the Enstore/SAM "infrastructure", the CDF and D0 approaches to data handling differ markedly, reflecting very different approaches to physics analysis.

D0 plans a 10kbyte/event "thumbnail" dataset that must be disk resident and should lessen the need to access the ~150kbyte/event DST datasets that will normally be on tape only and is expected to be accessed by continuous cycling through a small disk staging area. This minimizes disk costs, but puts a very high load on the relatively small (5% of the data) disk cache and on the tape system. At this stage the impact of heavy load on the disk system has not been modeled.

CDF plan that physics analysis should focus on the use of the 100 kbyte/event PAD format, all of which should be available on disk. Most of the access is expected to be to derived PAD files in which individual analysis groups will have concentrated the events needed for their work. Since the disk resident PAD datasets will be huge, reaching more than half a petabyte, there is less worry that heavy access to individual disk servers will cause performance problems. However, it does seem likely that much of physics analysis needing access to the PAD format will need much less than the full 100kbytes of information but will nevertheless have to read this from the disks and transport it across the network.

While it appears unlikely that two such different models of data access can both be optimum, it does not seem appropriate to force convergence to a single approach at this stage.

Both CDF and D0 assume that commodity disk servers at a current price of $10k for 1.9 terabytes will be a suitable technology. At this (decreasing with Moore's Law) price, the CDF model of putting all analysis datasets on disk is financially reasonable. However some cost increase may be needed to maintain adequate reliability and manageability.

## *Farms and Production*

### Production Farms

This section covers the production farms for both data reconstruction and Monte Carlo generation.

## Reconstruction Farms

The current design and operation of the production farms appears to be in good shape. Both CDF and D0 showed that the current production farms are capable of keeping up with the present data taking rates of the two experiments. They are to be congratulated on the achievement.

The current speed of the reconstruction code varies somewhat between the two experiments. CDF showed 2.5 GHz-sec/event while D0 is currently at 6.5 GHz-sec/event for the reconstruction-only part of the processing time (another 3.5 GHz-sec/event is spent on ntuple production and file merging making the total time 10 GHz-sec/event). In both experiments the processing time is dominated by the time spent during tracking. Given the different nature of the tracking detectors in the two experiments we do not expect the D0 reconstruction time to necessarily be as short as the CDF time. However we would hope that some inprovements could be made to keep the D0 execution time closer to the 10 GHz-sec/event rather than the assumed time of 25 GHz-sec/event. The estimated increase in execution time for D0 reconstruction of a factor of 4 with luminosity is worrisome and needs to be watched closely. We note that CDF assumed an increase in execution time of a factor of two at high luminosity for 132 nsec bunch spacing. We were not shown the scaling with luminosity for the CDF execution time. D0 used times based on 396 nsec bunch spacing. We encourage CDF to also develop projections assuming the higher number of interactions that are present at 396 nsec.

The reprocessing estimates for both experiments are similar. CDF quoted 0.3 and D0 0.5. We note that there appears to be non-negligible reprocessing activity of selected datasets being carried out on the analysis systems also that is not included in this number.

## Monte Carlo Farms

D0 plans to carry out all of its Monte Carlo production on remote farms located at collaborating institutions. They have already generated about 18M events through their full Geant simulation. This would appear to be an effective use of offsite resources and D0 are to be commended for this. We encourage D0 to tune their parametrized simulation so that they can reach their goals of 25 Hz of Monte Carlo (25% full Geant and 75% parametrized) for Run IIa. We note that should the parametrized simulation not measure up then 4 times more CPU will be required to achieve the 25 Hz goal.

CDF plans to generate a similar amount of Monte Carlo. The plan presented assumed that all of the Monte Carlo will be generated on the FNAL reconstruction farm. We were not shown any plans for generating MC at remote institutions and we encourage CDF to explore such use of remote resources.

## *Central Analysis Facilities*

CDF and DZero presented their plans for central analysis systems used to provide the storage and CPU for post-reconstruction physics analysis. Both experiments derived the specifications using input from their physics groups and run I experience. The data analyses in the physics groups are just beginning, and, as speakers from both experiments pointed out, this results in significant uncertainty in the assumptions used to specify the systems. Thus, it will be important to revisit the designs after a year of physics-quality data taking and analysis.

Both experiments intend to physically separate the disk and CPU for the systems, using Linux-based, commodity PC batch farms for the majority of the CPU. The committee supports this model. The disk model is quite different for the two experiments. DZero has found the reliability and I/O performance of the existing SMP machine to be sufficient for their high-data-volume analyses through roughly 2005. CDF has had difficulty with their similar SMP system, and instead of relying of this for the long run, they intend to deploy commodity PC fileservers, each

with 2 TB of directly connected disk and gigabit ethernet through which the disk is served to the batch CPUs.

Initial tests indicate this is feasible, but it is not yet proven in general use. The committee recommends that CDF continue to pursue this approach, and if proven to work, it is recommended that DZero consider this as a replacement of the SMP system (as mentioned in their presentation).

The DZero plan has two Linux-based batch clusters of analysis CPU: the SMP back end CAB, and the back end to the desktop cluster, CluB. The latter is to be managed by the experiment, providing a means to use the CPU and commodity disk provided by institutions. It will be used to analyze the data sets of roughly 1 TB needed by physics groups. The committee feels that this is an artificial split and the commodity-based analysis cluster for D0 should be a centrally managed system.

The DZero plan requires roughly $500k/year from collaborators for analysis systems. The committee feels that this should not broken into many small, uncoordinated purchases, but should be a smaller number of coordinated acquisitions to insure compatibility. This would be particularly important if CAB and CluB were merged.

Finally, the committee recommends that the experiments and the Computing Division consider all costs, including maintenance and personnel, when comparing various system options.

### *Remote Analysis and Grid Development*

The offline computing plans for the Dzero and CDF experiments recognize the importance of integrating plans for offline data analysis facilities at Fermilab with the remote computing facilities available to their collaborators. Dzero, in particular, has identified the role of "Global Systems" in its management structure, but both collaborations expected that international colleagues will play a significant role in the creation of Monte Carlo data sets, post-processing of selected event samples, and physics analysis using these remote computing resources.

Dzero has had in place large-scale MC event generation and simulation at 6 remote sites, producing since November 2001 over 17.8M fully simulated events. The committee congratulates the Dzero on this milestone and on its aggressive strategy to develop Regional Analysis Centres (RAC) that would provide centralized regional access to data analysis resources. CDF's effort has been more modest, making a commitment to the prototype deployment of Dzero's SAM as the remote interface to datasets located at Fermilab, and has made significant progress on this project over the last six months. Both collaborations have developed the tools to successfully maintain up-to-date analysis environments at a large number of home institutions, an important step especially during a period of rapid software development.

The committee is impressed with the support both collaborations have given to the integration of these off-site resources into the overall analysis plans. However, many of the commitments made to this activity appear to be informal, and not necessarily part of an institution's contribution to the collaboration. The committee recommends that both collaborations develop more detailed plans for the coordinated use of remote computing facilities, and identify how best they can be used to meet the overall Run II data processing and physics analysis needs. These plans should be formalized through Memoranda of Understanding (MOU's) and should identify clearly what remote resources can be made generally available to the entire collaboration. This is particularly important in the context of Dzero, where the very large estimates for physics analysis computing resources exceed the overall resources readily available to the experiment at Fermilab.

The integration of the remote computing resources into the Run II physics analysis effort will require significant improvements to the wide area network (WAN) connectivity of the laboratory. The Committee understood that Fermilab's connection to ESNET will quadruple in speed to an

OC12 link (622 Mbps) in the next several months.  However, the connectivity of ESNET to most of the collaborating institutions is hampered by the lack of a high-speed link from ESNET to STARLIGHT in Chicago, which then provides high-speed WAN Internet II and CANET access to North American institutions as well as Europe.  The Committee strongly endorses the plans to upgrade ESNET's capability, and recommends the Laboratory to work with ESNET to develop a high-speed connection to STARLIGHT on the timescale of the next 6 months.

Development of GRID-enabled analysis tools for CDF and Dzero is at an early stage,  appropriate for such an R&D effort.  Small groups in CDF, Dzero and Computing Division are collaborating in other GRID R&D efforts, with the immediate focus being the transformation of SAM into a GRID-enabled software tool. As a more functional set of GRID middleware becomes available, the two experiments appear to appropriately positioned to contribute effectively to the development effort and provide early examples of working implementations.  The Committee believes that the GRID is a promising technology for managing and using remotely distributed computing resources, and encourages both collaborations to continue their respective efforts and their close collaboration. The first step of GRID-enabling SAM is a reasonable goal and is endorsed by the Committee.

### *Infrastructure*
Overall the plans presented for network upgrades seem reasonable and appropriate.
The upgrade of networks in the trailers and assembly buildings to a modern switched network with 100Mbit connections for physicist desktops is clearly justified, however the more aggressive ideas to upgrade to Gigabit desktop network connections seems unreasonable.  In addition the requirements on the network from a CLUB facility not located in FCC may be more expensive than the potential benefit of such a cluster would warrant (compare costs of network infrastructure with cost of cluster).

### *Budgets, Staffing, Schedule, Planning*

### Budgets
In looking at the budgets presented by the experiments it was clear that there is a very large uncertainty in what the exact computing requirements will be and that those requirements will surely change as the experiments accumulate data and the collaborations get more experience and understanding.  Given the assumptions made by the experiments, which are not unreasonable starting points, the $2M equipment budgets seem reasonable at least for the first year or two.  However, we would expect that as the groups begin to analyse significantly more data over that time, their computing models will evolve and will require that these estimates be adjusted.  We would expect to see the first indications of that in a second review one year from now.

From the information presented it seemed that the CDF budget estimates were close to the guideline values derived from their initial assumptions, whereas for D0 the estimates were somewhat larger and they have explicitly assigned some $500K per year as costs to be borne by external institutions.  It is possible that this asymmetry is simply due to the differences between the two detectors and the resulting longer reconstruction and analysis times for D0.  The D0 simulation is also very expensive and outside of this budget envelope.

Since it is likely that contributions to the computing from external institutions will become important to the success of the analysis, we recommend that the collaborations make this contribution somewhat more formal, perhaps in the form of MOU's as discussed above.  We understand that moves in this direction are under discussion.

We considered whether the costs of computing depend on integrated luminosity or simply on running time.  In their presentations, the experiments had assumed scaling with luminosity.  This is not an entirely obvious assumption since both groups will run at a constant DAQ rate no matter

what the luminosity.  However, in discussions, it became clear that there are other effects, for example, as the data sample grow s, more and more analyses become feasible and more computing power and storage is needed to support that.  Again, given the large uncertainties here, the assumption of scaling with integrated luminosity seems realistic until experience shows otherwise.

We would also recommend considering some mechanism whereby the full costs of hardware and software choices are recognized by the experiments.  The issue of the maintenance costs of the large SMP's are a good example of costs hidden to the experiments.  Would they make the same choices if such costs had to be recognized by them?

## Staffing
- CDF should be encouraged to contribute dedicated effort to the SAM project, since now both experiments are making use of it.
- CD is providing some 35 FTE of effort dedicated to each experiment.  It is not entirely obvious that these levels are optimum.  Several historical considerations have led to the situation where potential synergies and consequent savings are missed.  Good examples would be common farms and analysis clus ters, where although they may be rigidly split into D0 and CDF pieces the management would be simplified – coordinated hardware purchases, etc.  Successful example is RHIC Computing Facility where a single cluster is used by all 4 experiments for both reconstruction and analysis.  Significant efficiency gains are possible.
- With limited staff, it is unreasonable to expect CD support for uncoordinated hardware purchases – desktops, small clusters, etc.  As well as being a significant security risk (and thus s taff cost), it is not unreasonable to put certain minimum requirements on systems and their management.

## Schedules and Planning
Given the initial assumptions and their uncertainties, the schedules for equipment and tape purchases seem reasonable.  However, as mentioned above, we would expect to see the plans refined continuously over the next few years.  The experiments are developing models for estimating these resource needs, however there are very large uncertainties in the assumptions made and revisiting this a year from now in the light of experience gained in that time seems necessary.